

Adversarial Machine Learning in Cybersecurity: Attacks and Defenses

Hu Ke¹, Jian Xu², Yong Wang³, Heyao Chen⁴, Zepeng Shen⁵

¹Mechanical Design, Manufacturing and Automation, Heilongjiang Institute of Technology, Heilongjiang, China

²Electrical and Electronics Engineering, University of Southern California, California, USA

³Information Technology, University of Aberdeen, Aberdeen, United Kingdom

⁴Computer Science and Technology, Beijing University of Posts and Telecommunications, Beijing, China

⁵Network Engineering, Shaanxi University of Technology, Shaanxi 723001, China

Abstract: *Adversarial Machine Learning (AML) refers to the research field that involves testing and improving machine learning models by introducing adversarial samples or attack techniques. In the cybersecurity domain, AML has significant potential to help identify and defend against threats such as malware, cyber attacks, and identity fraud. However, AML also faces numerous challenges, including low efficiency in generating adversarial samples, insufficient stealth, and issues with the generality and adaptability of defense methods. There is a dynamic interplay between adversarial attacks and defenses, with attackers continually developing new techniques and defenders needing to constantly improve their defense strategies. This interaction drives the rapid development of AML technology, making it increasingly important in cybersecurity. By deeply studying the interplay between adversarial attacks and defenses, the robustness and reliability of cybersecurity systems can be effectively enhanced, laying the foundation for future AI development in cybersecurity.*

Keywords: Adversarial Machine Learning, Cybersecurity, Robustness.

1. INTRODUCTION

1.1 Background Introduction

In recent years, the application of Artificial Intelligence (AI) and Machine Learning (ML) technologies in the field of cybersecurity has become increasingly widespread. As the complexity and frequency of cyber attacks continue to rise, traditional security defense methods are gradually revealing their inadequacies. AI and ML technologies can significantly improve the protective capabilities of cybersecurity systems by automating the analysis of large amounts of data, identifying potential threats, and predicting attack patterns. For example, ML-based malware detection systems can efficiently classify unknown samples by learning from known malware characteristics; identity authentication systems can enhance the accuracy and security of user authentication using biometric techniques; and network traffic analysis systems can promptly detect and respond to cyber attacks through real-time monitoring and anomaly detection. However, with the extensive application of AI and ML technologies, Adversarial Machine Learning (AML) has gradually become an important topic in the cybersecurity field. AML refers to the practice where attackers design adversarial samples to cause ML models to make errors in prediction and classification. Adversarial attacks exploit the vulnerabilities of ML models by introducing small but targeted perturbations to input data, leading to significant changes in model predictions. These attacks not only threaten the effectiveness of ML-based cybersecurity systems but can also result in severe security vulnerabilities, allowing attackers to bypass defense mechanisms and achieve their malicious objectives [1].

1.2 Problem Definition

The core concept of adversarial attacks is the **adversarial sample**, which is a new sample generated by adding small, imperceptible perturbations to a normal sample. These adversarial samples can cause an otherwise accurate ML model to make errors in prediction, thereby achieving the attacker's specific goals. Adversarial attacks can be categorized into various types based on the attacker's knowledge of the target model and the attack strategy. For instance, attacks can be classified as white-box or black-box based on the attacker's level of knowledge about the target model, and as evasion, poisoning, or model extraction attacks based on the attack strategy. The research objective of this paper is to explore the application of adversarial machine learning in the cybersecurity domain, analyze its potential threats, and propose effective defense methods. Specifically, the paper will provide a detailed analysis of the application of adversarial attacks in scenarios such as malware detection, identity authentication systems, and network traffic classification, discuss the advantages and disadvantages of different defense techniques, and validate the effectiveness of these defense methods through experiments and evaluations. By

analyzing the interplay between attacks and defenses, this paper aims to provide valuable references and guidance for researchers and practitioners in the cybersecurity field [2].

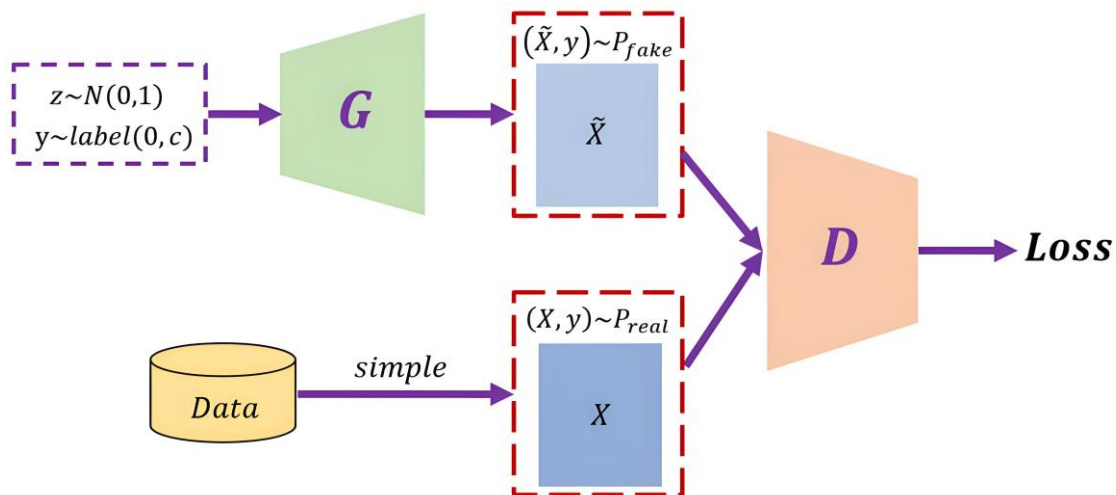
1.3 Paper Structure

This paper first introduces the background and importance of adversarial machine learning in cybersecurity. It then defines the core concepts of adversarial attacks and states the research objectives and overall structure of the paper. Specifically, the paper delves into the various application scenarios of adversarial attacks in cybersecurity, analyzes effective defense methods, and evaluates the practical effectiveness of these methods through experiments. The aim is to provide readers with a comprehensive perspective on the application of adversarial machine learning in the cybersecurity domain [3].

2. FOUNDATIONS OF ADVERSARIAL MACHINE LEARNING

2.1 Definition and Generation of Adversarial Samples

An adversarial sample is a new sample generated by adding small, imperceptible perturbations to a normal sample. These perturbations, though negligible to humans, can cause significant changes in the predictions of machine learning models. This phenomenon reveals the vulnerability of machine learning models, which are sensitive to small changes in high-dimensional input spaces. Common methods for generating adversarial samples include the Fast Gradient Sign Method (FGSM), Projected Gradient Descent (PGD), Carlini & Wagner (C&W) attack, and DeepFool. FGSM generates adversarial samples by computing the gradient of the loss function and adding perturbations along the gradient direction. It is simple and efficient but introduces larger perturbations. PGD optimizes the perturbations through multiple iterations, making the generated adversarial samples more subtle and effective. The C&W attack is a high-precision adversarial attack method that optimizes the loss function to generate adversarial samples that can successfully deceive the model, but it is computationally expensive. DeepFool incrementally increases the perturbations until the model misclassifies the sample, offering a good balance between performance and computational efficiency. These generation algorithms have their own characteristics but share the common goal of generating adversarial samples that can deceive machine learning models [4].



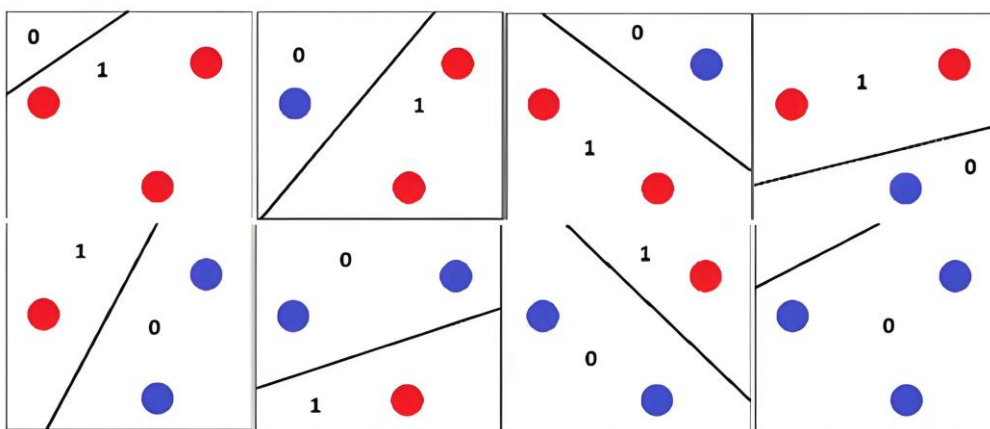
2.2 Classification of Adversarial Attacks

Adversarial attacks can be classified based on the attacker's knowledge of the target model and the attack strategy. From the perspective of target-driven attacks, adversarial attacks can be divided into white-box and black-box attacks. A white-box attack assumes that the attacker has complete knowledge of the target model's structure, parameters, and training data, allowing them to generate targeted adversarial samples using this information. This type of attack is typically more effective but requires extensive prior knowledge. A black-box attack assumes that the attacker has no knowledge of the target model's internal structure and can only generate adversarial samples by querying the model's input-output relationship. Black-box attacks are more challenging but are more common in practical applications because attackers typically cannot obtain complete information about the model. From the

perspective of strategy-driven attacks, adversarial attacks can be categorized into evasion attacks, poisoning attacks, and model extraction attacks. An evasion attack involves generating adversarial samples to cause the model to misclassify during the inference phase, primarily targeting deployed models. A poisoning attack involves polluting the training data during the training phase to make the model learn incorrect features, leading to poor performance during inference. A model extraction attack involves reconstructing a similar model by querying the target model and then using this model to generate adversarial samples. These three attack strategies have their own characteristics but share the common goal of causing machine learning models to make errors through different means [5].

$$H = \{set\ of\ linear\ classifiers\ in\ 2\ Dimensions\}$$

Then any 3 points can be classified by H correctly with separating hyperplane as shown in the following figure.



2.3 Characteristics of Adversarial Attacks in Cybersecurity

In the cybersecurity domain, adversarial attacks exhibit dynamic, diverse, and highly stealthy characteristics. Dynamic nature refers to the continuous development of adversarial attack methods and techniques, with new attack techniques emerging and attackers quickly adapting and improving their strategies. Diversity is reflected in the ability of attackers to design different attack methods for various security systems and application scenarios, such as malware detection systems, identity authentication systems, and network traffic classification systems. High stealth is achieved by generating adversarial samples that visually resemble normal samples, making them difficult to detect by humans or simple detection methods, thereby increasing the difficulty of defense. These characteristics make adversarial attacks particularly threatening in cybersecurity, and security systems need to be constantly updated and optimized to counteract the evolving attack techniques [6].

3. APPLICATIONS OF ADVERSARIAL ATTACKS IN CYBERSECURITY

3.1 Adversarial Attacks in Malware Detection

In the field of malware detection, machine learning models are extensively used in both static and dynamic analysis to identify and classify malware. However, adversarial attacks can render these models ineffective by generating adversarial samples, allowing malware to evade detection. Specifically, in static analysis, attackers can make minor modifications to the bytecode or signature of an executable to create adversarial samples that cause the model to misclassify the malicious nature of the file. In dynamic analysis, attackers can modify the behavior patterns of malware to mimic benign software, thereby deceiving behavior analysis models. A typical case involves using Generative Adversarial Networks (GANs) to generate adversarial malware. In this approach, the generator network learns to create malware samples that appear similar to legitimate software, while the discriminator network attempts to distinguish between genuine benign software and generated adversarial malware. Through this adversarial training process, the generated adversarial samples can successfully bypass traditional malware detection systems, demonstrating high evasion capabilities. Additionally, studies have shown

that modifying the code structure of malware or using obfuscation techniques can create adversarial samples. These techniques alter the surface features of malware, making it more similar to legitimate software at the bytecode level, thereby confusing machine learning-based detection models. This type of attack poses a significant threat to cybersecurity because it can render even the most advanced malware detection systems unable to detect hidden threats. Therefore, developing effective defense mechanisms to counter these advanced adversarial attacks is a critical area of research in cybersecurity [7].

3.2 Adversarial Attacks in Identity Authentication Systems

Identity authentication systems are crucial components in securing cybersecurity, with biometric technologies such as fingerprint and face recognition being widely adopted due to their uniqueness and convenience. However, these systems are also susceptible to adversarial attacks. Attackers can use adversarial samples to deceive authentication systems by adding minor perturbations to biometric features, thereby gaining unauthorized access. In fingerprint recognition, attackers can create fake fingerprints containing adversarial perturbations to deceive fingerprint sensors. For instance, researchers have demonstrated the ability to add specific texture patterns to silicone fake fingerprints, causing machine learning models to misclassify them as real user fingerprints. Similarly, in face recognition, attackers can generate adversarial images or videos to trick face recognition systems into identifying the attacker as an authorized user. This can be achieved by adding specific patterns to glasses or using special makeup on the face. These adversarial perturbations are usually imperceptible to humans but are sufficient to mislead machine learning models. Additionally, other biometric technologies such as voice and voiceprint recognition are also at risk of similar adversarial attacks. Attackers can introduce minor noise into voice samples to alter the model's recognition results. These attacks not only threaten the unlocking and access control of personal devices but also impact more critical security systems such as bank authentication and access control in government agencies. Ensuring the security of biometric systems, particularly their resilience against adversarial attacks, is an important research direction [8].

3.3 Other Scenarios

Beyond malware detection and identity authentication systems, adversarial attacks have significant impacts in other areas of cybersecurity. For example, in network traffic classification, machine learning models are used to identify different types of network traffic, including normal traffic and potential attack traffic. Attackers can generate adversarial network traffic, which is carefully designed to bypass traffic classification models without raising suspicion, thereby enabling various network attacks such as denial-of-service attacks and data theft. Additionally, threat intelligence data poisoning is another important application of adversarial attacks. Threat intelligence systems rely on accurate and reliable data to identify and respond to security threats. Attackers can inject false threat intelligence data to disrupt the system's normal operation, causing security teams to make incorrect decisions or expend resources to address false threats, thereby masking real attacks. This type of attack can render security response mechanisms ineffective, significantly impacting an organization's cybersecurity defense capabilities. In the Internet of Things (IoT) domain, the risk of adversarial attacks is more pronounced due to the limited resources and diverse characteristics of IoT devices. Attackers can tamper with sensor data to create adversarial inputs, causing the decision systems of devices to make incorrect judgments, thereby controlling the devices or triggering security incidents. In summary, adversarial attacks have multifaceted applications in cybersecurity and exhibit high levels of stealth and effectiveness. With the increasing application of machine learning in cybersecurity, effectively defending against these attacks to ensure system security is an urgent issue [9].

4. METHODS FOR DEFENDING AGAINST ADVERSARIAL ATTACKS

4.1 Classification of Defense Techniques

Defense methods against adversarial attacks can be categorized into two main types: robustness enhancement methods and detection and filtering methods. Robustness enhancement methods aim to improve the model's ability to resist adversarial attacks by enhancing the training process. Adversarial training is a common technique that involves introducing adversarial samples during training to improve the model's robustness. Specifically, adversarial training requires the model not only to learn the features of normal samples but also to correctly classify adversarial samples. This method effectively enhances the model's performance against adversarial attacks by exposing it to potential attack forms during training, thereby improving its generalization capabilities. Another robustness enhancement method is data augmentation, which improves the model's generalization by

expanding the training dataset. Data augmentation can be achieved by generating new training samples or transforming existing samples. For example, operations such as rotating, scaling, and cropping images can generate more training samples, making the model more stable when facing different input forms. This method not only enhances the model's robustness but also reduces the risk of overfitting, thereby improving the model's overall performance [10]. Detection and filtering methods focus on detecting and filtering adversarial samples during the model inference phase. Adversarial sample detection techniques analyze the features of input data to identify potential adversarial samples. For example, anomaly detection techniques can be used to identify adversarial samples by detecting abnormal patterns in the input data. Another method is input transformation, which involves preprocessing the input data to eliminate or reduce the impact of adversarial perturbations before it enters the model. Techniques such as image denoising and data normalization can reduce the impact of adversarial samples on the model. While these methods can improve model security to some extent, they often require additional computational resources and time [11].

4.2 Specific Defense Practices in Cybersecurity

In the field of cybersecurity, specific defense practices need to be customized according to different application scenarios. In malware detection, defense strategies can include multi-layer detection mechanisms and dynamic analysis techniques. Multi-layer detection mechanisms combine static and dynamic analysis to improve detection accuracy and robustness. Static analysis can detect the code features of malware, while dynamic analysis monitors the runtime behavior of malware to identify potential threats. Additionally, techniques such as adversarial training and data augmentation can improve the robustness of malware detection models, making them better able to withstand adversarial attacks. In identity authentication systems, security enhancement schemes can include multi-factor authentication and the use of diverse biometric features. Multi-factor authentication combines passwords, biometric features, and other authentication methods to enhance system security. For example, in face recognition systems, combining it with fingerprint or voice recognition can increase the difficulty of attacks. Furthermore, adversarial training and input transformation techniques can improve the robustness of biometric systems, enabling them to better resist adversarial attacks. For instance, in face recognition systems, image preprocessing techniques can eliminate the impact of adversarial perturbations, thereby improving recognition accuracy [12].

4.3 Game Theory Analysis of Defense and Attack

The dynamic process of defense measures and attack methods is a continuous evolution. Attackers continuously develop new attack methods to bypass existing defense measures, while defenders refine their defense technologies to counter new attack forms. This game-theoretic process drives the rapid development of adversarial machine learning, with both defense technologies and attack methods continually advancing. For example, as adversarial training techniques become more widespread, attackers may develop new attack methods to bypass these defenses, while defenders need to continuously update and optimize their defense strategies to counter new threats [13]. In the game between defense and attack, cost-effectiveness analysis is a critical consideration. Implementing defense measures often requires additional computational resources and time, while attackers need to invest time and effort to develop new attack methods. Therefore, defenders must strike a balance between enhancing defense effectiveness and reducing resource consumption. For example, while adversarial training can effectively improve model robustness, its training process is usually complex and time-consuming, requiring a balance between cost and benefit in practical applications. Similarly, while input transformation techniques can reduce the impact of adversarial samples, the preprocessing process may increase system latency, thereby affecting user experience. Therefore, when designing and implementing defense measures, various factors need to be considered to achieve optimal defense effectiveness [14].

5. EXPERIMENTS AND EVALUATION

5.1 Dataset and Experimental Setup

Choosing appropriate datasets and models is crucial when evaluating the effectiveness of adversarial attack and defense methods. In this study, the cybersecurity datasets used include malware detection datasets, network traffic classification datasets, and biometric datasets. For malware detection, we utilized public datasets such as ANDROZOO and VirusShare, which contain a large number of real malware samples and benign applications. In network traffic classification, we employed datasets like KDD Cup 1999 and CICIDS2017, which include various types of network traffic, including normal traffic and various attack traffic [15]. In the field of biometrics, we used

image datasets such as MNIST and CIFAR-10, as well as publicly available fingerprint and face recognition datasets, such as FVC2000 and LFW. In terms of experimental setup, we selected various commonly used machine learning models, including Convolutional Neural Networks (CNN), Recurrent Neural Networks (RNN), and Support Vector Machines (SVM). These models have demonstrated good performance in different cybersecurity tasks. The experimental environment utilized standard machine learning frameworks such as TensorFlow and PyTorch, and the experiments were conducted on high-performance computing servers. The evaluation metrics mainly include Accuracy, Precision, Recall, F1 Score, and robustness indicators such as the attack success rate of adversarial samples and the effectiveness of defense methods. Through comprehensive analysis of these metrics, the performance of adversarial attacks and defense methods can be thoroughly evaluated. [16]

5.2 Analysis of Attack Effectiveness

To evaluate the effectiveness of adversarial attacks, we conducted attack experiments on different cybersecurity models. First, in malware detection models, we used common adversarial attack methods such as FGSM, PGD, and C&W to generate adversarial samples. The experimental results showed that these adversarial samples could significantly reduce the detection accuracy of the models. For example, adversarial samples generated using FGSM reduced the accuracy of the CNN model from 95% to 70%, while samples generated using PGD further reduced the accuracy to 50%. This indicates that even simple adversarial attack methods can pose a serious threat to malware detection models. In network traffic classification, we also used the aforementioned adversarial attack methods to generate adversarial network traffic [17]. The experimental results showed that adversarial samples had a significant impact on the performance of network traffic classification models. In particular, the classification accuracy of the model decreased by more than 30% when using adversarial samples generated by PGD. This demonstrates that adversarial attacks can not only bypass malware detection systems but also deceive network traffic classification models, thereby enabling various network attacks. In the field of biometrics, we used adversarial attack methods to generate images to test the robustness of face recognition and fingerprint recognition systems. The results showed that adversarial samples could significantly reduce the recognition rate of these systems, thereby enabling unauthorized access [18].

5.3 Evaluation of Defense Effectiveness

To evaluate the effectiveness of various defense methods, we conducted experiments on adversarial training, data augmentation, adversarial sample detection, and input transformation. First, we performed adversarial training on malware detection models. The experimental results showed that models trained with adversarial training exhibited significant robustness improvements when facing adversarial samples. For example, the attack success rate of models against adversarial samples generated by FGSM decreased from 70% to 30%, while the success rate against samples generated by PGD decreased from 50% to 20%. This indicates that adversarial training can effectively enhance the defense capabilities of models. In terms of data augmentation, we expanded the existing malware detection dataset to improve the generalization ability of the model. The experimental results showed that the performance degradation of the model when facing adversarial samples was significantly reduced. For example, the attack success rate of the model against adversarial samples generated by FGSM decreased from 70% to 40%. This demonstrates that data augmentation, by increasing the diversity of training samples, can improve the model's robustness against new attacks. In adversarial sample detection, we used feature analysis and anomaly detection techniques to identify adversarial samples. The experimental results showed that these methods could detect adversarial samples to some extent, thereby reducing their impact on model performance [19]. For example, the feature analysis method could detect approximately 80% of adversarial samples, which could then be filtered out. Additionally, input transformation methods, which preprocess input data, could significantly reduce the impact of adversarial perturbations. For example, through image denoising and normalization, the attack success rate of adversarial samples generated by FGSM could be reduced from 70% to 40%. When evaluating the effectiveness of various defense methods, it is also necessary to consider the trade-off between robustness and computational overhead. Although adversarial training and data augmentation can significantly improve the robustness of models, the training process of these methods is complex and time-consuming. For example, adversarial training requires multiple iterations of generating and training adversarial samples, leading to a significant increase in training time. Similarly, data augmentation methods require additional data processing steps, increasing computational resource consumption. Therefore, in practical applications, it is necessary to comprehensively consider the effectiveness and computational overhead of defense methods based on specific task requirements to find the best defense strategy [20].

6. CHALLENGES AND FUTURE TRENDS

6.1 Ongoing Challenges

Currently, adversarial attacks and defenses in the field of cybersecurity face several ongoing challenges. First, the generation efficiency and stealthiness of adversarial samples are pressing issues that need to be addressed. Existing adversarial sample generation methods, such as FGSM and PGD, can generate effective adversarial samples to some extent, but their generation efficiency is low, and they are easily detected and defended against. Future research needs to develop more efficient and stealthy adversarial sample generation methods that can successfully bypass security detection systems without raising suspicion. Additionally, the generalizability and adaptability of defense methods are also significant challenges. Existing defense methods are often optimized for specific types of attacks and models, lacking generality and being difficult to adapt to different scenarios and tasks. Future research needs to develop more generalizable and adaptable defense methods to address diverse adversarial attacks.

6.2 Future Development Directions

In the future, the development of adversarial machine learning in cybersecurity will mainly focus on the following directions. First, research on more efficient adversarial training methods to enhance model robustness will be a key focus. Although existing adversarial training methods can improve the defense capabilities of models, their training process is complex and time-consuming. Future research needs to develop more efficient adversarial training methods that can train highly robust models in a shorter time. Second, the defense against adversarial attacks in federated learning will become an important research direction. Federated learning, as a distributed machine learning method, protects user privacy while also facing threats from adversarial attacks. Future research needs to explore how to effectively detect and defend against adversarial attacks in federated learning to ensure system security. Finally, the ethical and regulatory norms for adversarial machine learning will also become an important trend in the future. As adversarial attack and defense technologies continue to develop, ensuring that these technologies are applied in accordance with ethical standards and relevant laws and regulations will be a key topic for researchers and policymakers. By establishing reasonable norms and standards, the healthy development of adversarial machine learning technology in cybersecurity can be guided.

7. CONCLUSION

Adversarial attacks pose a severe threat in cybersecurity, and existing research has revealed their potential risks in areas such as malware detection, network traffic classification, and biometrics. Through the evaluation of various types of defense measures, such as adversarial training, data augmentation, adversarial sample detection, and input transformation, we found that these methods have significant effects in improving model robustness, but challenges such as low generation efficiency, poor stealthiness, and insufficient generality still exist. The importance of adversarial machine learning cannot be overlooked, as it not only relates to the security of networks but also affects the healthy development of AI technology. Future research should focus on developing more efficient adversarial training methods, exploring defense strategies in federated learning, and establishing corresponding ethical and regulatory norms to promote the further development of cybersecurity AI.

REFERENCES

- [1] Yu Pengwen. Legal nature and application rules of artificial intelligence evidence in criminal proceedings [J]. Chinese Journal of Criminal Law, 2024, (05): 36-54. DOI: 10.19430/j.cnki.3891.2024.05.010.
- [2] Wei Kuo, et al. Strategic application of AI intelligent algorithm in network threat detection and defense [J]. Journal of Theory and Practice of Engineering Science, 2024, 4(01): 49-57.
- [3] Chen Wangmei, et al. Applying machine learning algorithm to optimize personalized education recommendation system [J]. Journal of Theory and Practice of Engineering Science, 2024, 4(01): 101-108.
- [4] Sun Jin. Application of machine learning in network anomaly detection [J]. Information and Computer (Theory Edition), 2024, 36(09): 81-83.
- [5] Hao Ning. Research on network security management application based on artificial intelligence technology [J]. Information Recording Materials, 2024, 25(02): 66-68. DOI: 10.16009/j.cnki.cn13-1295/tq.2024.02.021.
- [6] Wang Lingtong, Wang Huiling, Xu Miao, et al. Overview of detection and defense technologies for cross-site scripting attacks [J]. Journal of Computer Applications, 2024, 41(03): 652-662. DOI: 10.19734/j.issn.1001-3695.2023.06.0286.

- [7] Tian Miao, et al. The application of artificial intelligence in medical diagnostics: A new frontier [J]. *Academic Journal of Science and Technology*, 2023, 8(2): 57-61.
- [8] Chen Heyao, et al. Threat detection driven by artificial intelligence: Enhancing cybersecurity with machine learning algorithms [J]. 2024.
- [9] Huang Qimeng, Wu Miaomiao, Li Yun. Research on filtering adversarial feature selection for evasion attacks [J]. *Telecommunications Science*, 2023, 39(07): 46-58.
- [10] Gao Ying, Chen Xiaofeng, Zhang Yiyu, et al. Overview of attack and defense technologies in federated learning systems [J]. *Journal of Computer Science*, 2023, 46(09): 1781-1805.
- [11] Du Shuqian, et al. Improving science question ranking with model and retrieval-augmented generation [C]. *The 6th International Scientific and Practical Conference "Old and New Technologies of Learning Development in Modern Conditions,"* Berlin, Germany: International Science Group, 2024: 252.
- [12] Cheng Shiwei, et al. 3D Pop-Ups: Omnidirectional image visual saliency prediction based on crowdsourced eye-tracking data in VR [J]. *Displays*, 2024, 83: 102746. Elsevier.
- [13] Hu Shengqiu, Li Youguo, Gao Yuan, et al. Internet of Things security detection method based on adversarial deep learning [J]. *Electronic Design Engineering*, 2022, 30(11): 50-54+59. DOI: 10.14022/j.issn1674-6236.2022.11.011.
- [14] Lin Sifang, et al. Artificial intelligence and electroencephalogram analysis: Innovative methods for optimizing anesthesia depth [J]. *Journal of Theory and Practice in Engineering and Technology*, 2024, 1(4): 1-10.
- [15] Zhou Siming, Li Dan. Attacks and defenses in public clouds based on machine learning [J]. *Network Security Technology & Application*, 2022, (01): 70-72.
- [16] Liu Qixu, Wang Junnan, Yin Jie, et al. Application of adversarial machine learning in network intrusion detection [J]. *Journal of Communications*, 2021, 42(11): 1-12.
- [17] Du Shuqian, et al. Improving science question ranking with model and retrieval-augmented generation [C]. *The 6th International Scientific and Practical Conference "Old and New Technologies of Learning Development in Modern Conditions,"* Berlin, Germany: International Science Group, 2024: 252.
- [18] Zhao Lemen, Wang Rui. Data science platform: Features, technologies, and trends [J]. *Computer Science*, 2021, 48(08): 1-12.
- [19] Cheng Shiwei, et al. Poster graphic design with your eyes: An approach to automatic textual layout design based on visual perception [J]. *Displays*, 2023, 79: 102458. Elsevier.
- [20] Yang Qi, Jia Peng, Liu Jiayong. Generation of DOCX adversarial samples based on DCGAN [J]. *Modern Computer*, 2021, (15): 77-81.